
Generative Adversarial Privacy

Chong Huang¹ Peter Kairouz² Xiao Chen² Lalitha Sankar¹ Ram Rajagopal²

Abstract

We present a data-driven framework called generative adversarial privacy (GAP). Inspired by recent advancements in generative adversarial networks (GANs), GAP allows the data holder to learn the privatization mechanism directly from the data. Under GAP, finding the optimal privacy mechanism is formulated as a constrained minimax game between a privatizer and an adversary. We show that for appropriately chosen adversarial loss functions, GAP provides privacy guarantees against strong information-theoretic adversaries. We also evaluate GAP’s performance on the GENKI face database.

1. Introduction

The use of machine learning algorithms for data analytics has recently seen unprecedented success for a variety of problems of practical relevance such as image classification, natural language processing, and prediction of consumer behavior, electricity use, political preferences, to name a few. The success of these algorithms hinges on the availability of large datasets, which are often crowd-sourced and contain private information. This, in turn, has led to privacy concerns and a growing body of research focused on developing privacy-guaranteed learning techniques.

Moving towards randomization-based methods, in recent years, two distinct approaches with provable *statistical privacy* guarantees have emerged: (a) context-free approaches that assume worst-case dataset statistics and adversaries; (b) context-aware approaches that explicitly model the dataset statistics and adversary’s capabilities. On the one hand, context-free approaches (such as differential privacy (Dwork & Roth, 2014)) provide strong privacy guarantees against worst-case adversaries, but often lead to a significant reduction in the utility and increased sample complexity (Fienberg et al., 2010; Wang et al., 2015; Yu et al., 2014; Karwa & Slavković, 2016; Duchi et al., 2016; Kairouz et al., 2016).

¹Arizona State University, Tempe, AZ, USA ²Stanford University, Stanford, CA, USA. Correspondence to: Lalitha Sankar <lalithasankar@asu.edu>.

On the other, context-aware approaches (such as mutual information privacy (Rebollo-Monedero et al., 2010; Calmon & Fawaz, 2012; Sankar et al., 2013)) achieve a better privacy-utility tradeoff by incorporating the statistics of the dataset and explicitly modeling the public and private variables, but necessitate knowledge of data statistics (such as joint priors over the public and private variables). Such information is hardly ever present in practice.

Given the challenges of existing approaches, we take a fundamentally new approach towards enabling private data publishing. Instead of adopting worst-case, context-free notions of data privacy, we introduce a novel context-aware model of privacy that allows the designer to cleverly add noise where it matters. We overcome the issue of statistical knowledge by taking a *data-driven approach*; specifically, we leverage recent advancements in generative adversarial networks (GANs) (Goodfellow et al., 2014; Mirza & Osindero, 2014) to introduce a framework for context-aware privacy that we call *generative adversarial privacy* (GAP).

2. Generative Adversarial Privacy

We consider a dataset \mathcal{D} which contains both public and private variables for n individuals. We represent the public variables by a random variable X , and the private variables (which are typically correlated with the public variables) by a random variable Y . Each dataset entry contains a pair of public and private variables denoted by (X, Y) . We assume that each entry pair (X, Y) is distributed according to $P(X, Y)$, and is independent from other entry pairs in the dataset. We define the privacy mechanism as a randomized mapping given by $\hat{X} = g(X, Y)$.

We define $\hat{Y} = h(g(X, Y))$ to be the adversary’s inference of the private variable Y from \hat{X} using a decision rule h . We allow for *hard decision rules* under which $h(g(X, Y))$ is a direct estimate of Y and *soft decision rules* under which $h(g(X, Y)) = P_h(\cdot|g(X, Y))$ is a distribution over \mathcal{Y} . To quantify the adversary’s performance, we use a loss function $\ell(h(g(X = x)), Y = y)$ defined for every public-private pair (x, y) . Thus, the expected loss of the adversary with respect to (*w.r.t.*) X and Y is

$$L(h, g) \triangleq \mathbb{E}[\ell(h(g(X, Y)), Y)], \quad (1)$$

where the expectation is taken over $P(X, Y)$ and the randomness in g and h .

The data holder would like to find a privacy mechanism g that is both privacy preserving (in the sense that it is difficult for the adversary to learn Y from \hat{X}) and utility preserving (in the sense that it does not distort the original data too much). In contrast, for a fixed choice of privacy mechanism g , the adversary would like to find a (potentially randomized) function h that minimizes its expected loss, which is equivalent to maximizing the negative of the expected loss. This leads to a constrained minimax game between the privatizer and the adversary given by

$$\begin{aligned} \min_{g(\cdot)} \max_{h(\cdot)} -L(h, g) \\ \text{s.t. } \mathbb{E}[d(g(X, Y), X)] \leq D, \end{aligned} \quad (2)$$

where the constant $D \geq 0$ determines the allowable distortion for the privatizer and the expectation is taken over $P(X, Y)$ and the randomness in g and h .

Theorem 1. *Under the class of hard decision rules, when $\ell(h(g(x, y), y))$ is the 0-1 loss function, the GAP minimax problem in (2) simplifies to*

$$\begin{aligned} \min_{g(\cdot)} \max_{y \in \mathcal{Y}} P(y, g(X, Y)) \\ \text{s.t. } \mathbb{E}[d(g(X, Y), X)] \leq D, \end{aligned} \quad (3)$$

indicating that maximizing the probability of correctly guessing Y is the optimal adversarial strategy for any privatizer, i.e., the adversary uses the MAP decision rule. On the other hand, for a soft-decision decoding adversary (i.e., $h = P_h(y|\hat{x})$ is a distribution over \mathcal{Y}) under log-loss function $\ell(h(g(X, Y), y)) = \log \frac{1}{P_h(y|g(X, Y))}$, the optimal adversarial strategy h^* is the posterior belief of Y given $g(X, Y)$ and the GAP minimax problem in (2) is equivalent to

$$\begin{aligned} \min_{g(\cdot)} I(g(X, Y); Y) \\ \text{s.t. } \mathbb{E}[d(g(X, Y), X)] \leq D, \end{aligned} \quad (4)$$

where $I(g(X, Y); Y)$ is the mutual information (MI) between $g(X, Y)$ and Y .

The above theorem shows that GAP can recover MI privacy (under a log loss) and MAP privacy (under a 0-1 loss). The proof of Theorem 1 is omitted for brevity.

2.1. Data-driven GAP

In the absence of $P(X, Y)$, we propose a data-driven version of GAP that allows the data holder to learn privatization mechanisms directly from a dataset $\mathcal{D} = \{(x_{(i)}, y_{(i)})\}_{i=1}^n$. Under the data-driven version of GAP, we represent the privacy mechanism via a generative model $g(X, Y; \theta_p)$ parameterized by θ_p . In the training phase, the data holder learns the optimal parameters θ_p by competing against a *computational adversary*: a classifier modeled by a neural network $h(g(X, Y; \theta_p); \theta_a)$ parameterized by θ_a .

In the data-driven approach, we can quantify the adversary's

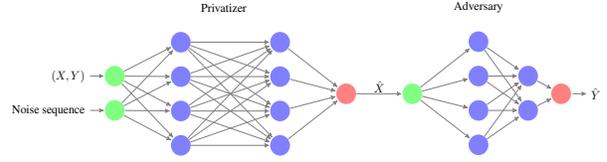


Figure 1: A multi-layer neural network model for the privatizer and adversary

empirical loss by

$$L_n(\theta_p, \theta_a) = -\frac{1}{n} \sum_{i=1}^n \ell(h(g(x_{(i)}, y_{(i)}; \theta_p); \theta_a), y_{(i)}) \quad (5)$$

where $(x_{(i)}, y_{(i)})$ is the i^{th} row of \mathcal{D} . The optimal parameters for the privatizer and adversary are the solutions to

$$\begin{aligned} \min_{\theta_p} \max_{\theta_a} -L_n(\theta_p, \theta_a) \\ \text{s.t. } \mathbb{E}_{\mathcal{D}}[d(g(X, Y; \theta_p), X)] \leq D, \end{aligned} \quad (6)$$

where the expectation is over \mathcal{D} and the randomness in g .

3. GAP for the GENKI Dataset

To demonstrate GAP's capability of learning the privacy mechanism directly from the data, we test our model on the GENKI dataset (Whitehill & Movellan, 2012) which contains 1940 grey-scale images of faces. We consider gender as private variable Y and the image pixels as public variable X . Two different privatizer architectures are studied: the feedforward neural network privatizer (FNNP) and the transposed convolutional neural network privatizer (TCNNP). The FNNP uses a five-layer feedforward neural network to combine the low-dimensional noise with the original image. The TCNNP uses a three-layer transposed convolutional neural network to generate high-dimensional additive noise from low-dimensional noise. The adversary is modeled by a seven-layer convolutional neural network.

Figure 2 illustrates the gender classification accuracy of the adversary for different values of distortion. We observe that the adversary's accuracy of classifying the private label (gender) decreases progressively as the distortion increases. Given the same distortion value, FNNP achieves better privacy protection compared with TCNNP. The adversary's miss-classified privatized image samples are shown in Figure 3. We observe that both privatizers change mostly eyes, nose, mouth, beard, and hair.

References

Calmon, F. P. and Fawaz, N. Privacy against statistical inference. In *Communication, Control, and Computing (Allerton), 2012 50th Annual Allerton Conference on*, pp. 1401–1408, 2012.

Duchi, John, Wainwright, Martin, and Jordan, Michael. Min-

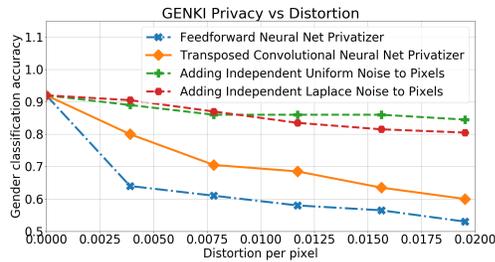


Figure 2: Privacy-distortion tradeoff

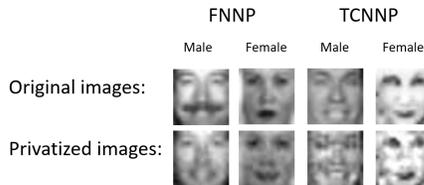


Figure 3: Miss-classified privatized image samples

imax optimal procedures for locally private estimation. *arXiv preprint arXiv:1604.02390*, 2016.

Dwork, Cynthia and Roth, Aaron. The algorithmic foundations of differential privacy. *Found. Trends Theor. Comput. Sci.*, 9(3–4):211–407, August 2014. ISSN 1551-305X. doi: 10.1561/04000000042. URL <http://dx.doi.org/10.1561/04000000042>.

Fienberg, Stephen E., Rinaldo, Alessandro, and Yang, Xiaolin. *Differential Privacy and the Risk-Utility Trade-off for Multi-dimensional Contingency Tables*, pp. 187–199. Springer Berlin Heidelberg, Berlin, Heidelberg, 2010. ISBN 978-3-642-15838-4. doi: 10.1007/978-3-642-15838-4_17. URL https://doi.org/10.1007/978-3-642-15838-4_17.

Goodfellow, Ian, Pouget-Abadie, Jean, Mirza, Mehdi, Xu, Bing, Warde-Farley, David, Ozair, Sherjil, Courville, Aaron, and Bengio, Yoshua. Generative adversarial nets. In *Advances in neural information processing systems*, pp. 2672–2680, 2014.

Kairouz, Peter, Bonawitz, Keith, and Ramage, Daniel. Discrete distribution estimation under local privacy. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48, ICML’16*, pp. 2436–2444. JMLR.org, 2016. URL <http://dl.acm.org/citation.cfm?id=3045390.3045647>.

Karwa, Vishesh and Slavković, Aleksandra. Inference using noisy degrees: Differentially private β -model and synthetic graphs. *The Annals of Statistics*, 44(1):87–112, 2016.

Mirza, Mehdi and Osindero, Simon. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.

Rebollo-Monedero, D., Forne, J., and Domingo-Ferrer, J. From t-Closeness-Like Privacy to Postrandomization via Information Theory. *IEEE Transactions on Knowledge and Data Engineering*, 22(11):1623–1636, November 2010. ISSN 1041-4347. doi: 10.1109/TKDE.2009.190.

Sankar, L., Rajagopalan, S. R., and Poor, H. V. Utility-privacy tradeoffs in databases: An information-theoretic approach. *IEEE Transactions on Information Forensics and Security*, 8(6):838–852, 2013.

Wang, Yue, Lee, Jaewoo, and Kifer, Daniel. Differentially private hypothesis testing, revisited. *arXiv preprint arXiv:1511.03376*, 2015.

Whitehill, Jacob and Movellan, Javier. Discriminately decreasing discriminability with learned image filters. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pp. 2488–2495. IEEE, 2012.

Yu, Fei, Fienberg, Stephen E, Slavković, Aleksandra B, and Uhler, Caroline. Scalable privacy-preserving data sharing methodology for genome-wide association studies. *Journal of biomedical informatics*, 50:133–141, 2014.